

1

This study analyzes chronic kidney disease using machine learning techniques based on a chronic kidney disease (CKD).

Chronic Kidney Disease (CKD) or chronic renal disease has become a major issue with a steady growth rate. A person can only survive without kidneys for an average time of 18 days, which makes a huge demand for a kidney transplant and Dialysis. It is important to have effective methods for early prediction of CKD. Machine learning methods are effective in CKD prediction. This work proposes a workflow to predict CKD status based on clinical data, incorporating data preprocessing, a missing value handling method with collaborative filtering and attributes selection

1

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).

1

2

jupyter nbextension install --py luxwidget jupyter nbextension enable --py luxwidget pip install lux-api

```
1 # import necessary packages
2
3 import os
4 import sys
5 import numpy as np
6 import pandas as pd
7 import matplotlib.pyplot as plt
8 import seaborn as sns
9 sns.set()
10 %matplotlib inline
11 import sklearn
12 import warnings
13 warnings.filterwarnings('ignore')
14 #import lux # EDA
15 from sklearn.preprocessing import MinMaxScaler

1 # import the dataset
2 df = pd.read_csv('/content/drive/MyDrive/Colab_Notebooks/DL/Kidney Disease Health Care/kidney_disease.csv')
3 df.head(n=10)
```

	id	age	bp	sg	al	su	rbc	pc	pcc	b
0	0	48.0	80.0	1.020	1.0	0.0	NaN	normal	notpresent	notpresen
1	1	7.0	50.0	1.020	4.0	0.0	NaN	normal	notpresent	notpresen
2	2	62.0	80.0	1.010	2.0	3.0	normal	normal	notpresent	notpresen

```

1 df.shape
(400, 26)
5 5 60.0 80.0 1.015 2.0 0.0 NaN NaN notpresent notpresen
1 # Total count of null values
2 df.isnull().sum()

id          0
age         9
bp         12
sg         47
al         46
su         49
rbc        152
pc         65
pcc         4
ba          4
bgr        44
bu         19
sc         17
sod        87
pot        88
hemo       52
pcv        70
wc        105
rc        130
htn         2
dm          2
cad         2
appet       1
pe          1
ane         1
classification 0
dtype: int64

1 # Total Null count in percentage
2 df.isnull().sum() / len(df) *100

id          0.00
age         2.25
bp          3.00
sg         11.75
al         11.50
su         12.25
rbc        38.00
pc         16.25
pcc         1.00
ba          1.00
bgr        11.00
bu          4.75
sc          4.25
sod        21.75
pot        22.00
hemo       13.00
pcv        17.50
wc         26.25
rc         32.50
htn         0.50
dm          0.50
cad         0.50
appet       0.25
pe          0.25
ane         0.25

```

```
classification    0.00
dtype: float64
```

```
1 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 400 entries, 0 to 399
Data columns (total 26 columns):
 #   Column                Non-Null Count  Dtype
---  -
 0   id                    400 non-null    int64
 1   age                  391 non-null    float64
 2   bp                   388 non-null    float64
 3   sg                   353 non-null    float64
 4   al                   354 non-null    float64
 5   su                   351 non-null    float64
 6   rbc                  248 non-null    object
 7   pc                   335 non-null    object
 8   pcc                  396 non-null    object
 9   ba                   396 non-null    object
10  bgr                  356 non-null    float64
11  bu                   381 non-null    float64
12  sc                   383 non-null    float64
13  sod                  313 non-null    float64
14  pot                  312 non-null    float64
15  hemo                 348 non-null    float64
16  pcv                  330 non-null    object
17  wc                   295 non-null    object
18  rc                   270 non-null    object
19  htn                  398 non-null    object
20  dm                   398 non-null    object
21  cad                  398 non-null    object
22  appet                399 non-null    object
23  pe                   399 non-null    object
24  ane                  399 non-null    object
25  classification       400 non-null    object
dtypes: float64(11), int64(1), object(14)
memory usage: 81.4+ KB
```

▼ Handling missing values using SimpleImputer

```
1 # sklearn approach to handle all missing data at one go
2
3 from sklearn.impute import SimpleImputer
4 imp_mode = SimpleImputer(missing_values=np.nan, strategy='most_frequent')
5 df_imputed = pd.DataFrame(imp_mode.fit_transform(df))
6 df_imputed.columns = df.columns
7 df_imputed
```

	id	age	bp	sg	al	su	rbc	pc	pcc	ba	...	pcv	wc	rc	htn	dm	cad	appe
0	0	48.0	80.0	1.02	1.0	0.0	normal	normal	notpresent	notpresent	...	44	7800	5.2	yes	yes	no	goc
1	1	7.0	50.0	1.02	4.0	0.0	normal	normal	notpresent	notpresent	...	38	6000	5.2	no	no	no	goc
2	2	62.0	80.0	1.01	2.0	3.0	normal	normal	notpresent	notpresent	...	31	7500	5.2	no	yes	no	po

▼ Imputed columns

```
...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...      ...
1 df_imputed.columns

Index(['id', 'age', 'bp', 'sg', 'al', 'su', 'rbc', 'pc', 'pcc', 'ba', 'bgr',
      'bu', 'sc', 'sod', 'pot', 'hemo', 'pcv', 'wc', 'rc', 'htn', 'dm', 'cad',
      'appet', 'pe', 'ane', 'classification'],
      dtype='object')

1 # After handling missing data, there is no null values
2 df_imputed.isnull().sum()

id          0
age         0
bp          0
sg          0
al          0
su          0
rbc         0
pc          0
pcc         0
ba          0
bgr         0
bu          0
sc          0
sod         0
pot         0
hemo        0
pcv         0
wc          0
rc          0
htn         0
dm          0
cad         0
appet       0
pe          0
ane         0
classification 0
dtype: int64

1 df_imputed.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 400 entries, 0 to 399
Data columns (total 26 columns):
#   Column          Non-Null Count  Dtype
---  -
0   id              400 non-null    object
1   age             400 non-null    object
2   bp              400 non-null    object
3   sg              400 non-null    object
4   al              400 non-null    object
5   su              400 non-null    object
6   rbc             400 non-null    object
7   pc              400 non-null    object
8   pcc             400 non-null    object
9   ba              400 non-null    object
10  bgr             400 non-null    object
11  bu              400 non-null    object
12  sc              400 non-null    object
13  sod             400 non-null    object
```

```
14 pot          400 non-null object
15 hemo         400 non-null object
16 pcv          400 non-null object
17 wc          400 non-null object
18 rc          400 non-null object
19 htn          400 non-null object
20 dm          400 non-null object
21 cad          400 non-null object
22 appet        400 non-null object
23 pe          400 non-null object
24 ane          400 non-null object
25 classification 400 non-null object
dtypes: object(26)
memory usage: 81.4+ KB
```

```
1 df.describe()
```

	id	age	bp	sg	al	su	bgr	bu	sc
count	400.000000	391.000000	388.000000	353.000000	354.000000	351.000000	356.000000	381.000000	383.000000
mean	199.500000	51.483376	76.469072	1.017408	1.016949	0.450142	148.036517	57.425722	3.072454
std	115.614301	17.169714	13.683637	0.005717	1.352679	1.099191	79.281714	50.503006	5.741126
min	0.000000	2.000000	50.000000	1.005000	0.000000	0.000000	22.000000	1.500000	0.400000
25%	99.750000	42.000000	70.000000	1.010000	0.000000	0.000000	99.000000	27.000000	0.900000
50%	199.500000	55.000000	80.000000	1.020000	0.000000	0.000000	121.000000	42.000000	1.300000
75%	299.250000	64.500000	80.000000	1.020000	2.000000	0.000000	163.000000	66.000000	2.800000
max	399.000000	90.000000	180.000000	1.025000	5.000000	5.000000	490.000000	391.000000	76.000000



Unique values in each columns

```
1 # Finding the unique values in the columns
2
3 for i in df_imputed.columns:
4     print("*****",i,"*****")
5     print()
6     print(set(df_imputed[i].tolist()))
7     print()
```



```

##### WC #####
{'12300', '5700', '6400', '21600', '6900', '9200', '11300', '\t?', '4200', '12700', '12400', '12500', '10200',
##### RC #####
{'4.6', '3', '5.2', '8.0', '6.4', '3.1', '2.9', '3.3', '6.2', '4.0', '6.0', '\t?', '2.1', '2.3', '3.5', '2.8',
##### htn #####
{'yes', 'no'}

##### dm #####
{' yes', '\tno', 'yes', '\tyes', 'no'}

##### cad #####
{'yes', '\tno', 'no'}

##### appet #####
{'good', 'poor'}

##### pe #####
{'yes', 'no'}

##### ane #####
{'yes', 'no'}

##### classification #####
{'ckd\t', 'notckd', 'ckd'}

```

Abbreviations:

pcv - nwq
 wc - nwq
 rc - num with quot
 dm- char
 cad - char
 classification - char

```

1 print(df_imputed['pcv'].mode())
2 print(df_imputed['wc'].mode())
3 print(df_imputed['rc'].mode())

0    41.0
Name: pcv, dtype: float64
0    9800.0
Name: wc, dtype: float64
0     5.2
Name: rc, dtype: float64

1 df_imputed['classification'].value_counts()

ckd      248
notckd   150
ckd\t      2
Name: classification, dtype: int64

```

```
1 df_imputed['classification'] = df_imputed['classification'].apply(lambda x: 'ckd' if x == 'ckd\t' else x)
```

```
1 df_imputed['cad'].value_counts()
```

```
no      364
yes      34
\tno      2
Name: cad, dtype: int64
```

```
1 df_imputed['cad'] = df_imputed['cad'].apply(lambda x: 'no' if x == '\tno' else x)
```

```
1 df_imputed['dm'].value_counts()
```

```
no      260
yes     134
\tno      3
\tyes      2
yes        1
Name: dm, dtype: int64
```

```
1 df_imputed['dm'] = df_imputed['dm'].apply(lambda x: 'no' if x == '\tno' else x)
2 df_imputed['dm'] = df_imputed['dm'].apply(lambda x: 'yes' if x == '\tyes' else x)
3 df_imputed['dm'] = df_imputed['dm'].apply(lambda x: 'yes' if x == ' yes' else x)
```

```
1 df_imputed['dm'].value_counts()
```

```
2
```

```
3 #PCV = Packed Cell Volume
```

```
4
```

```
no      263
yes     137
Name: dm, dtype: int64
```

▼ Abbreviation

pcv - nwq

wc - nwq

rc - num with quot

dm- char

cad - char

classification - char

```
1 df_imputed['rc'].value_counts()
```

```
5.2      148
4.5       16
4.9       14
4.7       11
4.8       10
3.9       10
4.6        9
3.4        9
5.9        8
5.5        8
6.1        8
5.0        8
3.7        8
5.3        7
5.8        7
5.4        7
3.8        7
```

```

5.6      6
4.3      6
4.2      6
3.2      5
4.4      5
5.7      5
6.4      5
5.1      5
6.2      5
6.5      5
4.1      5
3.6      4
6.3      4
6.0      4
4.0      3
3.3      3
4        3
3.5      3
2.9      2
3.1      2
2.6      2
2.1      2
2.5      2
2.8      2
3.0      2
2.7      2
5        2
2.3      1
\t?      1
2.4      1
3        1
8.0      1
Name: rc, dtype: int64

```

```
1 df_imputed['rc'] = df_imputed['rc'].apply(lambda x: '5.2' if x=='\t?' else x)
```

```
1 df_imputed['wc'].value_counts()
```

```

9800      116
6700       10
9600        9
7200        9
9200        9
...
19100       1
\t?         1
12300       1
14900       1
12700       1
Name: wc, Length: 92, dtype: int64

```

```
'\t?' '\t6200' '\t8400'
```

```
9800
```

```

1 df_imputed['wc'] = df_imputed['wc'].apply(lambda x: '9800' if x=='\t?' else x)
2 df_imputed['wc'] = df_imputed['wc'].apply(lambda x: '6200' if x=='\t6200' else x)
3 df_imputed['wc'] = df_imputed['wc'].apply(lambda x: '8400' if x=='\t8400' else x)

```

```
1 df_imputed['pcv'].value_counts()
```

```

41      91
52      21
44      19
48      19
40      16
43      14
42      13
45      13

```



```

32      12
36      12
33      12
50      12
28      12
34      11
37      11
30      9
29      9
35      9
46      9
31      8
24      7
39      7
26      6
38      5
53      4
51      4
49      4
47      4
54      4
25      3
22      3
27      3
19      2
23      2
15      1
21      1
17      1
20      1
\t43    1
18      1
9       1
\t?     1
16      1
14      1

```

Name: pcv, dtype: int64

```

1 df_imputed['pcv'] = df_imputed['pcv'].apply(lambda x: '43' if x=='\t43' else x)
2 df_imputed['pcv'] = df_imputed['pcv'].apply(lambda x: '41' if x=='\t?' else x)

```

```

1 for i in df_imputed.columns:
2     print("*****",i,"*****")
3     print()
4     print(set(df_imputed[i].tolist()))
5     print()

```

```

##### rc #####
{'4.6', '3', '5.2', '8.0', '6.4', '3.1', '2.9', '3.3', '6.2', '4.0', '6.0', '2.1', '2.3', '3.5', '2.8', '4.9',
##### htn #####
{'yes', 'no'}

##### dm #####
{'yes', 'no'}

##### cad #####
{'yes', 'no'}

##### appet #####
{'good', 'poor'}

##### pe #####
{'yes', 'no'}

##### ane #####
{'yes', 'no'}

##### classification #####
{'notckd', 'ckd'}

```

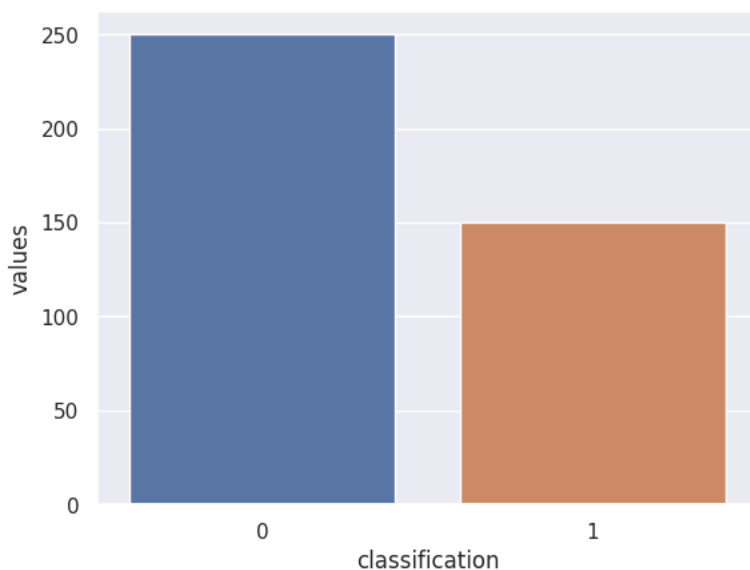
▼ Chronic Kidney Disease Vs Non Chronic Kidney Disease

```

1
2 temp = df_imputed['classification'].value_counts()
3 temp_df = pd.DataFrame({'classification': temp.index, 'values': temp.values})
4 print(sns.barplot(x = 'classification', y='values', data=temp_df))

```

Axes(0.125,0.11;0.775x0.77)



```
1 df_imputed['classification'].value_counts()
```

```

ckd      250
notckd   150
Name: classification, dtype: int64

```

```
1 df_imputed.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 400 entries, 0 to 399
Data columns (total 26 columns):
#   Column                Non-Null Count  Dtype
---  -
0   id                    400 non-null   object
1   age                  400 non-null   object
2   bp                   400 non-null   object
3   sg                   400 non-null   object
4   al                   400 non-null   object
5   su                   400 non-null   object
6   rbc                  400 non-null   object
7   pc                   400 non-null   object
8   pcc                  400 non-null   object
9   ba                   400 non-null   object
10  bgr                  400 non-null   object
11  bu                   400 non-null   object
12  sc                   400 non-null   object
13  sod                  400 non-null   object
14  pot                  400 non-null   object
15  hemo                 400 non-null   object
16  pcv                  400 non-null   object
17  wc                   400 non-null   object
18  rc                   400 non-null   object
19  htn                  400 non-null   object
20  dm                   400 non-null   object
21  cad                  400 non-null   object
22  appet               400 non-null   object
23  pe                   400 non-null   object
24  ane                  400 non-null   object
25  classification       400 non-null   object
dtypes: object(26)
memory usage: 81.4+ KB

```

```
1 df.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 400 entries, 0 to 399
Data columns (total 26 columns):
#   Column                Non-Null Count  Dtype
---  -
0   id                    400 non-null   int64
1   age                  391 non-null   float64
2   bp                   388 non-null   float64
3   sg                   353 non-null   float64
4   al                   354 non-null   float64
5   su                   351 non-null   float64
6   rbc                  248 non-null   object
7   pc                   335 non-null   object
8   pcc                  396 non-null   object
9   ba                   396 non-null   object
10  bgr                  356 non-null   float64
11  bu                   381 non-null   float64
12  sc                   383 non-null   float64
13  sod                  313 non-null   float64
14  pot                  312 non-null   float64
15  hemo                 348 non-null   float64
16  pcv                  330 non-null   object
17  wc                   295 non-null   object
18  rc                   270 non-null   object
19  htn                  398 non-null   object
20  dm                   398 non-null   object
21  cad                  398 non-null   object
22  appet               399 non-null   object
23  pe                   399 non-null   object
24  ane                  399 non-null   object

```

```

25 classification 400 non-null object
dtypes: float64(11), int64(1), object(14)
memory usage: 81.4+ KB

```

```
1 df.select_dtypes(exclude=["object"]).columns
```

```

Index(['id', 'age', 'bp', 'sg', 'al', 'su', 'bgr', 'bu', 'sc', 'sod', 'pot',
      'hemo'],
      dtype='object')

```

```

1 for i in df.select_dtypes(exclude=["object"]).columns :
2     df_imputed[i] = df_imputed[i].apply(lambda x:float(x))

```

```
1 df_imputed.dtypes
```

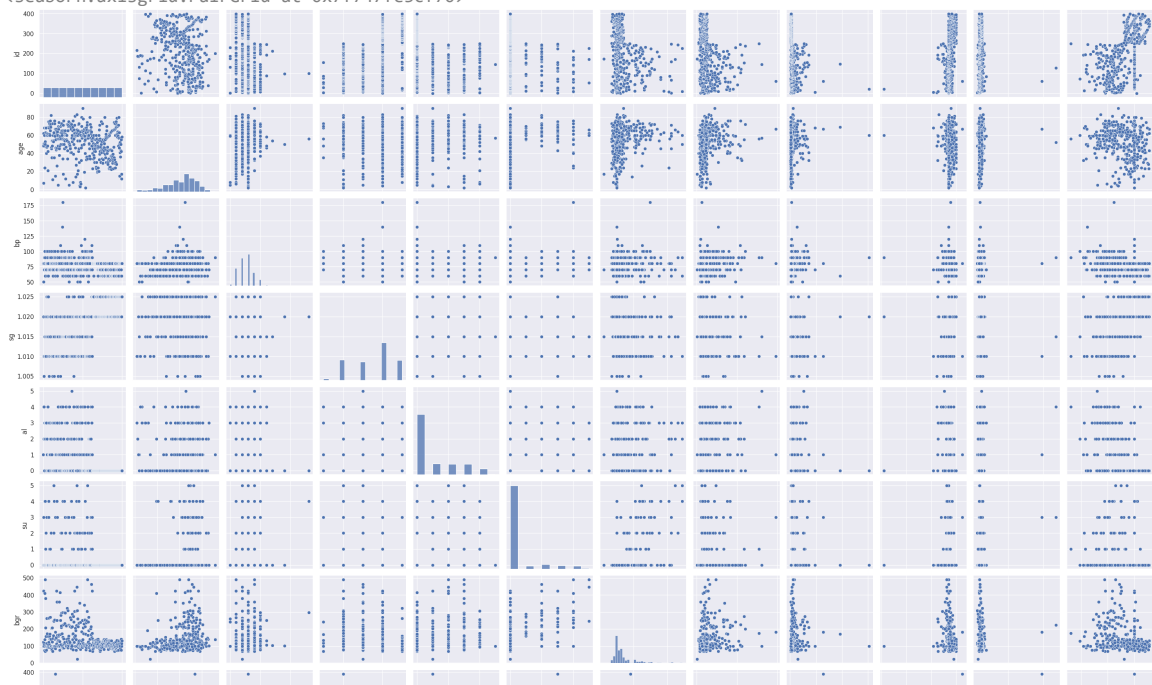
```

id                float64
age               float64
bp               float64
sg               float64
al               float64
su               float64
rbc              object
pc              object
pcc             object
ba              object
bgr             float64
bu             float64
sc             float64
sod             float64
pot             float64
hemo            float64
pcv             object
wc              object
rc              object
htn             object
dm              object
cad             object
appet           object
pe             object
ane            object
classification   object
dtype: object

```

```
1 sns.pairplot(df_imputed)
```

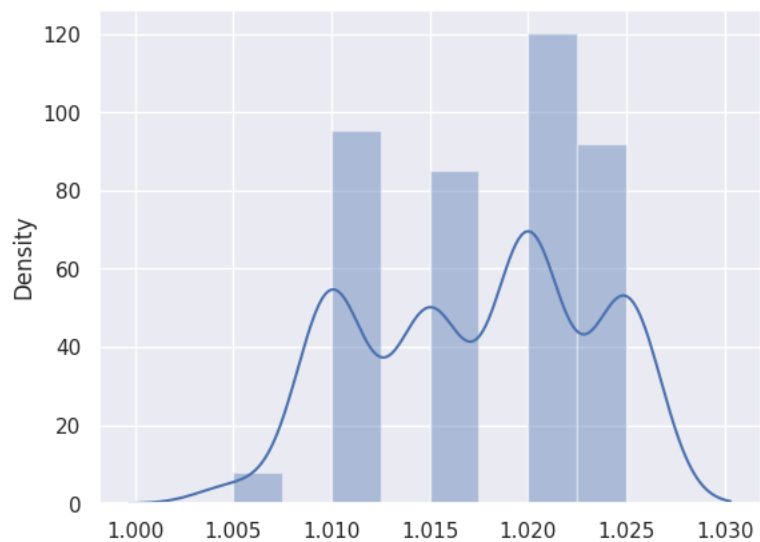
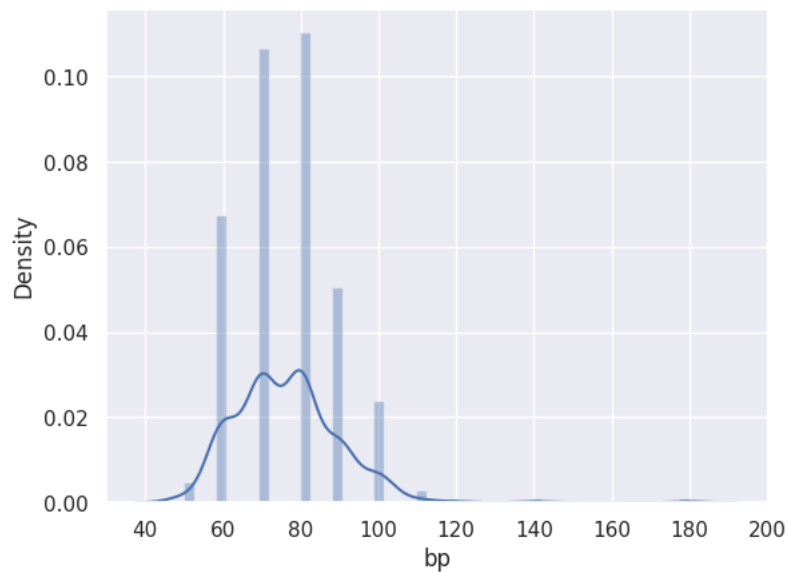
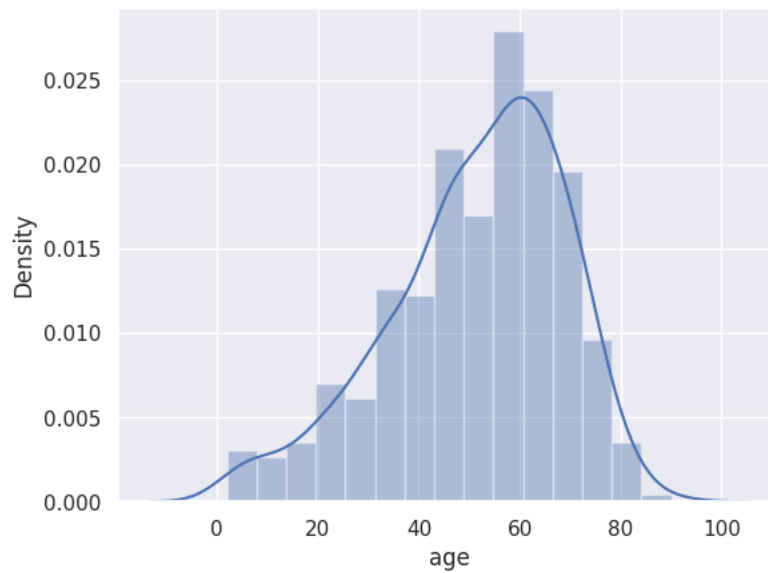
<seaborn.axisgrid.PairGrid at 0x7f747fe3cf70>



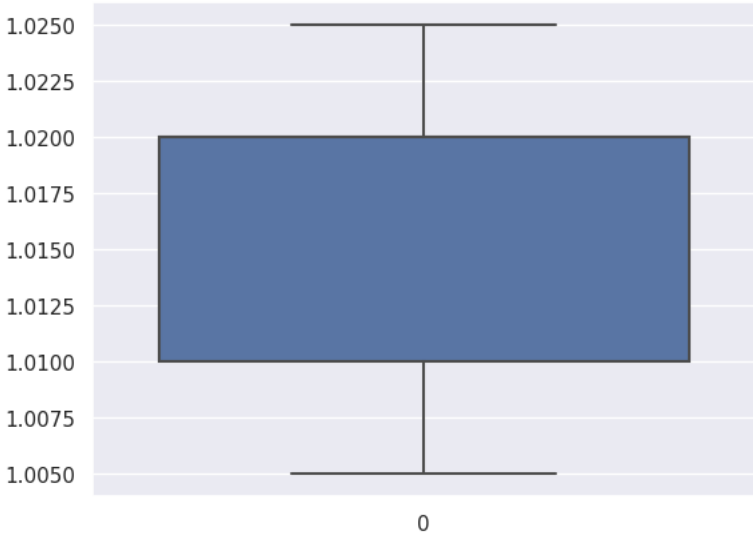
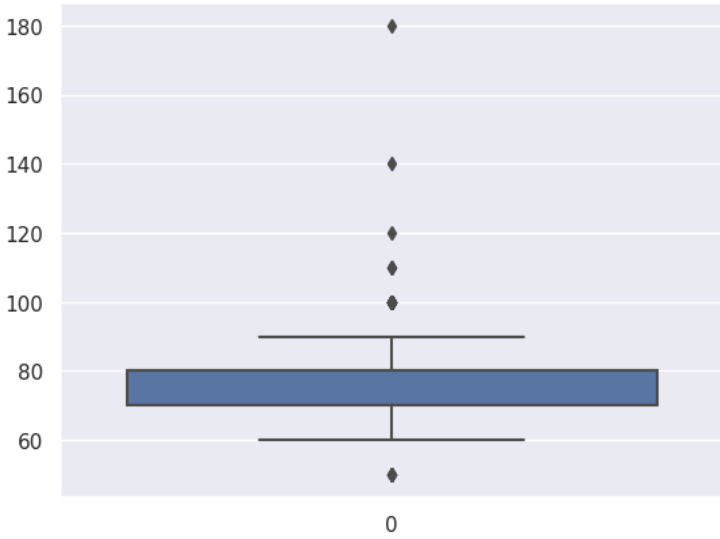
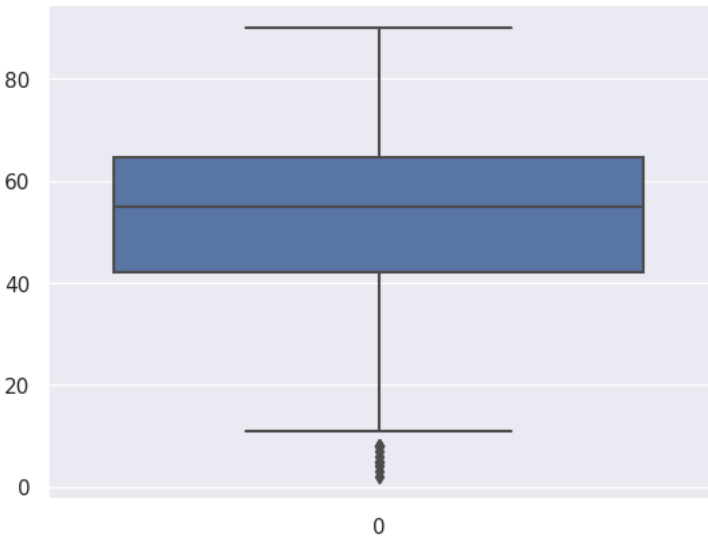
```

1 # Find the distribution of the data
2
3 def distplots(col):
4     sns.distplot(df[col])
5     plt.show()
6
7 for i in list(df_imputed.select_dtypes(exclude=['object']).columns)[1:]:
8     distplots(i)

```



```
1 # Find and remove outliers of data
2
3 def boxplots(col):
4     sns.boxplot(df[col])
5     plt.show()
6
7 for i in list(df_imputed.select_dtypes(exclude=['object']).columns)[1:]:
8     boxplots(i)
```




```
1 df_imputed.head(10)
```

	id	age	bp	sg	al	su	rbc	pc	pcc	ba	...	pcv	wc	rc	htn	dm	cad	ap
0	0.0	48.0	80.0	1.020	1.0	0.0	normal	normal	notpresent	notpresent	...	44	7800	5.2	yes	yes	no	gr
1	1.0	7.0	50.0	1.020	4.0	0.0	normal	normal	notpresent	notpresent	...	38	6000	5.2	no	no	no	gr
2	2.0	62.0	80.0	1.010	2.0	3.0	normal	normal	notpresent	notpresent	...	31	7500	5.2	no	yes	no	p
3	3.0	48.0	70.0	1.005	4.0	0.0	normal	abnormal	present	notpresent	...	32	6700	3.9	yes	no	no	p
4	4.0	51.0	80.0	1.010	2.0	0.0	normal	normal	notpresent	notpresent	...	35	7300	4.6	no	no	no	gr
5	5.0	60.0	90.0	1.015	3.0	0.0	normal	normal	notpresent	notpresent	...	39	7800	4.4	yes	yes	no	gr
6	6.0	68.0	70.0	1.010	0.0	0.0	normal	normal	notpresent	notpresent	...	36	9800	5.2	no	no	no	gr
7	7.0	24.0	80.0	1.015	2.0	4.0	normal	abnormal	notpresent	notpresent	...	44	6900	5	no	yes	no	gr
8	8.0	52.0	100.0	1.015	3.0	0.0	normal	abnormal	present	notpresent	...	33	9600	4.0	yes	yes	no	gr
9	9.0	53.0	90.0	1.020	2.0	0.0	abnormal	abnormal	present	notpresent	...	29	12100	3.7	yes	yes	no	p

10 rows x 26 columns



```
1 df_imputed.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 400 entries, 0 to 399
Data columns (total 26 columns):
#   Column              Non-Null Count  Dtype
---  -
0   id                   400 non-null    float64
1   age                  400 non-null    float64
2   bp                   400 non-null    float64
3   sg                   400 non-null    float64
4   al                   400 non-null    float64
5   su                   400 non-null    float64
6   rbc                  400 non-null    object
7   pc                   400 non-null    object
8   pcc                  400 non-null    object
9   ba                   400 non-null    object
10  bgr                  400 non-null    float64
11  bu                   400 non-null    float64
12  sc                   400 non-null    float64
13  sod                  400 non-null    float64
14  pot                  400 non-null    float64
15  hemo                 400 non-null    float64
16  pcv                  400 non-null    object
17  wc                   400 non-null    object
18  rc                   400 non-null    object
19  htn                  400 non-null    object
20  dm                   400 non-null    object
21  cad                  400 non-null    object
22  appet               400 non-null    object
23  pe                   400 non-null    object
24  ane                  400 non-null    object
25  classification       400 non-null    object
dtypes: float64(12), object(14)
memory usage: 81.4+ KB
```

▼ Label Encoding: To convert Categorical to Numerical

```
1 # Label encodint to convert categorical values to numerical
2
```

```
3 from sklearn import preprocessing
4 """
5 df_enco = df_imputed.apply(preprocessing.LabelEncoder().fit_transform)
6 df_enco
7 """
8
9 categorical_feature = [feature for feature in df_imputed.columns if df_imputed[feature].dtype == 'O']
10
```

```
1 df_imputed.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 400 entries, 0 to 399
Data columns (total 26 columns):
#   Column                Non-Null Count  Dtype
---  -
0   id                    400 non-null    float64
1   age                  400 non-null    float64
2   bp                   400 non-null    float64
3   sg                   400 non-null    float64
4   al                   400 non-null    float64
5   su                   400 non-null    float64
6   rbc                  400 non-null    object
7   pc                   400 non-null    object
8   pcc                  400 non-null    object
9   ba                   400 non-null    object
10  bgr                  400 non-null    float64
11  bu                   400 non-null    float64
12  sc                   400 non-null    float64
13  sod                  400 non-null    float64
14  pot                  400 non-null    float64
15  hemo                 400 non-null    float64
16  pcv                  400 non-null    object
17  wc                   400 non-null    object
18  rc                   400 non-null    object
19  htn                  400 non-null    object
20  dm                   400 non-null    object
21  cad                  400 non-null    object
22  appet               400 non-null    object
23  pe                   400 non-null    object
24  ane                  400 non-null    object
25  classification       400 non-null    object
dtypes: float64(12), object(14)
memory usage: 81.4+ KB
```

```
1 df_imputed['pcv'] = df_imputed['pcv'].astype('float')
2 df_imputed['wc'] = df_imputed['wc'].astype('float')
3 df_imputed['rc'] = df_imputed['rc'].astype('float')
```

```
1 df_imputed.head()
```

	id	age	bp	sg	al	su	rbc	pc	pcc	ba	...	pcv	wc	rc	htn	dm	cad	appe
0	0.0	48.0	80.0	1.020	1.0	0.0	normal	normal	notpresent	notpresent	...	44.0	7800.0	5.2	yes	yes	no	good
1	1.0	7.0	50.0	1.020	4.0	0.0	normal	normal	notpresent	notpresent	...	38.0	6000.0	5.2	no	no	no	good
2	2.0	62.0	80.0	1.010	2.0	3.0	normal	normal	notpresent	notpresent	...	31.0	7500.0	5.2	no	yes	no	poor
3	3.0	48.0	70.0	1.005	4.0	0.0	normal	abnormal	present	notpresent	...	32.0	6700.0	3.9	yes	no	no	poor
4	4.0	51.0	80.0	1.010	2.0	0.0	normal	normal	notpresent	notpresent	...	35.0	7300.0	4.6	no	no	no	good

5 rows × 26 columns

```
1 categorical_feature = [feature for feature in df_imputed.columns if df_imputed[feature].dtype == 'O']
2
```

```
1 categorical_feature
```

```
['rbc',
 'pc',
 'pcc',
 'ba',
 'htn',
 'dm',
 'cad',
 'appet',
 'pe',
 'ane',
 'classification']
```

```
18
```

```
1 df_imputed.info()
```

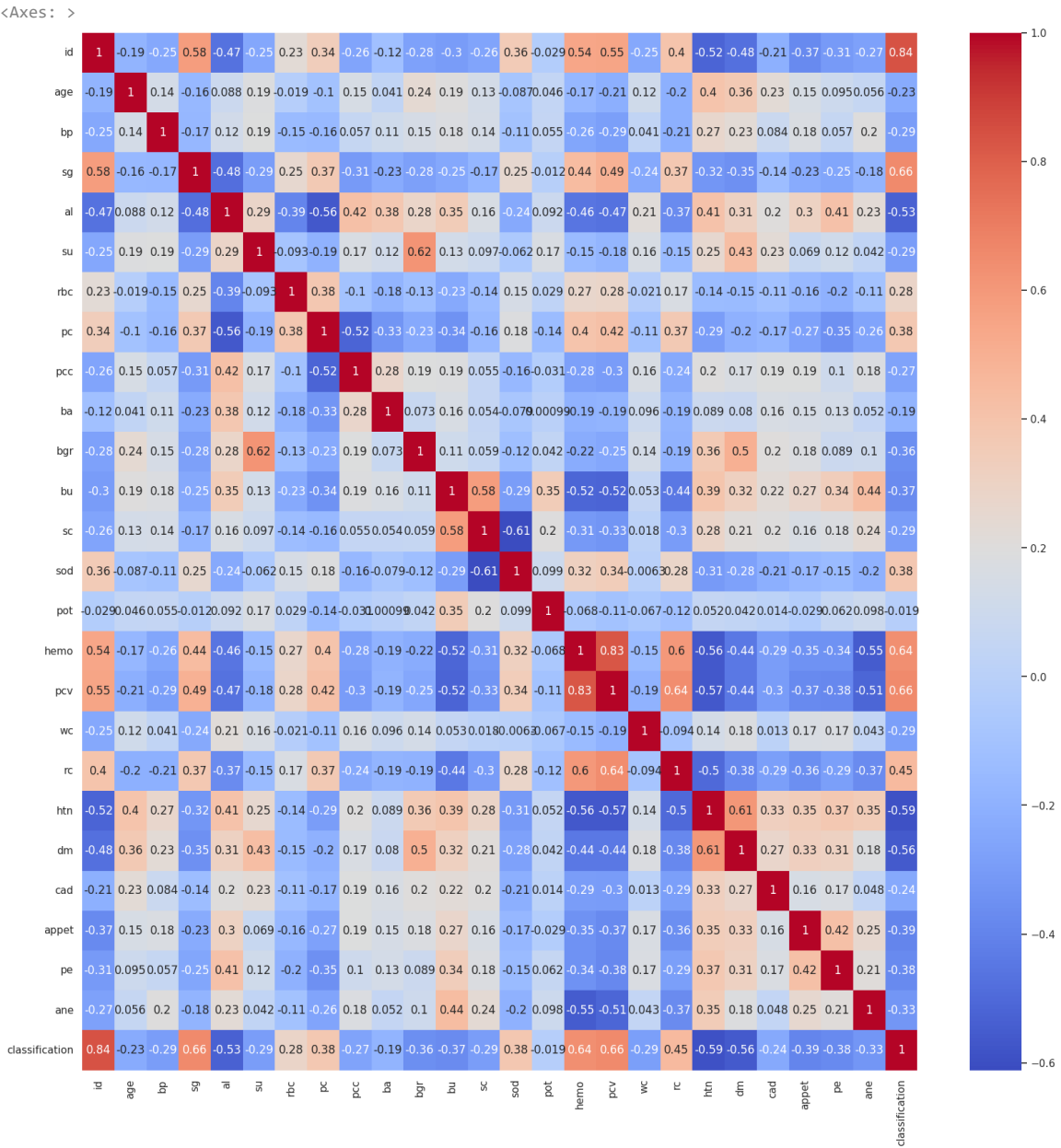
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 400 entries, 0 to 399
Data columns (total 26 columns):
#   Column                Non-Null Count  Dtype
---  -
0    id                    400 non-null   float64
1    age                   400 non-null   float64
2    bp                    400 non-null   float64
3    sg                    400 non-null   float64
4    al                    400 non-null   float64
5    su                    400 non-null   float64
6    rbc                   400 non-null   object
7    pc                    400 non-null   object
8    pcc                   400 non-null   object
9    ba                    400 non-null   object
10   bgr                   400 non-null   float64
11   bu                    400 non-null   float64
12   sc                    400 non-null   float64
13   sod                   400 non-null   float64
14   pot                   400 non-null   float64
15   hemo                  400 non-null   float64
16   pcv                   400 non-null   float64
17   wc                    400 non-null   float64
18   rc                    400 non-null   float64
19   htn                   400 non-null   object
20   dm                    400 non-null   object
21   cad                   400 non-null   object
22   appet                400 non-null   object
23   pe                   400 non-null   object
24   ane                   400 non-null   object
25   classification        400 non-null   object
dtypes: float64(15), object(11)
memory usage: 81.4+ KB
```

```
1 for i in df_imputed.select_dtypes(exclude=['float']).columns:
2     df_imputed[i] = df_imputed[i].astype('category')
3     df_imputed[i] = df_imputed[i].cat.codes
```

```
1 df_imputed.head()
```

	id	age	bp	sg	al	su	rbc	pc	pcc	ba	...	pcv	wc	rc	htn	dm	cad	appet	pe	ane	classification
0	0.0	48.0	80.0	1.020	1.0	0.0	1	1	0	0	...	44.0	7800.0	5.2	1	1	0	0	0	0	
1	1.0	7.0	50.0	1.020	1.0	0.0	1	1	0	0	...	38.0	6000.0	5.2	0	0	0	0	0	0	

```
1 # Finding the correlation
2
3 plt.figure(figsize=(20,20))
4 corr = df_imputed.corr()
5 sns.heatmap(corr, annot=True, cmap='coolwarm')
```



```
1 df_imputed['classification'].value_counts()
```

```

0    250
1    150
Name: classification, dtype: int64

1 # split the data into independent variable and dependent variable
2 x = df_imputed.drop(['id', 'classification'], axis=1)
3 y = df_imputed['classification']

1 x.head()

```

	age	bp	sg	al	su	rbc	pc	pcc	ba	bgr	...	hemo	pcv	wc	rc	htn	dm	cad	appet	pe	ane
0	48.0	80.0	1.020	1.0	0.0	1	1	0	0	121.0	...	15.4	44.0	7800.0	5.2	1	1	0	0	0	0
1	7.0	50.0	1.020	4.0	0.0	1	1	0	0	99.0	...	11.3	38.0	6000.0	5.2	0	0	0	0	0	0
2	62.0	80.0	1.010	2.0	3.0	1	1	0	0	423.0	...	9.6	31.0	7500.0	5.2	0	1	0	1	0	1
3	48.0	70.0	1.005	4.0	0.0	1	0	1	0	117.0	...	11.2	32.0	6700.0	3.9	1	0	0	1	1	1
4	51.0	80.0	1.010	2.0	0.0	1	1	0	0	106.0	...	11.6	35.0	7300.0	4.6	0	0	0	0	0	0

5 rows × 24 columns

```

1 y.head()

0    0
1    0
2    0
3    0
4    0
Name: classification, dtype: int8

```

▼ Label Balancing with RandomOverSampler

```

1 # lets detect the label balance
2
3 from imblearn.over_sampling import RandomOverSampler
4 from collections import Counter
5 print(Counter(y))

Counter({0: 250, 1: 150})

1 ros = RandomOverSampler()
2 x_ros , y_ros = ros.fit_resample(x, y)
3 print(Counter(y_ros))

Counter({0: 250, 1: 250})

```

▼ Feature Scaling using Standard Scaler

```

1 # feature scaling
2 # Min Max Scaler (-1,1)
3 # from sklearn.preprocessing import StandardScaler
4
5 scaler = MinMaxScaler((-1,1))
6 x = scaler.fit_transform(x_ros)
7 y = y_ros

1 x

```

```
array([[ 0.04545455, -0.53846154,  0.5      , ..., -1.      ,
        -1.      , -1.      ],
       [-0.88636364, -1.      ,  0.5      , ..., -1.      ,
        -1.      , -1.      ],
       [ 0.36363636, -0.53846154, -0.5      , ...,  1.      ,
        -1.      ,  1.      ],
       ...,
       [-0.54545455, -0.84615385,  1.      , ..., -1.      ,
        -1.      , -1.      ],
       [-0.40909091, -0.84615385,  1.      , ..., -1.      ,
        -1.      , -1.      ],
       [ 0.40909091, -0.84615385,  0.5      , ..., -1.      ,
        -1.      , -1.      ]])
```

```
1 y
```

```
0      0
1      0
2      0
3      0
4      0
```

```
..
```

```
495    1
496    1
497    1
498    1
499    1
```

```
Name: classification, Length: 500, dtype: int8
```

▼ This is all about the pre-processing approach

```
1 x.shape
```

```
(500, 24)
```

```
1 # Applying PCA method to reduce non-significant variables and values (if any)
2
3 from sklearn.decomposition import PCA
4
5 # it means the machine(sklearn) choose the minimum number of PCA such that 95% of the variance is retained
6
7 pca = PCA(0.95)
8 x_pca = pca.fit_transform(x)
9 print(x.shape)
10 print(x_pca.shape)

(500, 24)
(500, 15)
```

```
1 # split the data into training and test for building the model and prediction with test dataset
```

```
1 from sklearn.model_selection import train_test_split
2 x_train, x_test, y_train, y_test = train_test_split(x_pca, y, test_size=0.2, random_state=5)
```

Model Building

▼ MLP - multilayer perceptron (Neural Network)

```

1 import keras
2 from keras.models import Sequential
3 from keras.layers import Dense
4 from keras.layers import Dropout
5 from keras.callbacks import ModelCheckpoint, EarlyStopping
6 from keras.models import Model
7 from keras.optimizers import Adam

```

▼ creating the model

```

1 def model():
2     classifier = Sequential()
3     classifier.add(Dense(15, input_shape=(x_train.shape[1],), activation='relu'))
4     classifier.add(Dropout(0.2))
5     classifier.add(Dense(15, activation='relu'))
6     classifier.add(Dropout(0.4))
7     classifier.add(Dense(1, activation='sigmoid'))
8     classifier.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
9
10    return classifier
11

```

```

1 model = model()
2 model.summary()

```

Model: "sequential_2"

Layer (type)	Output Shape	Param #
=====		
dense_6 (Dense)	(None, 15)	240
dropout_4 (Dropout)	(None, 15)	0
dense_7 (Dense)	(None, 15)	240
dropout_5 (Dropout)	(None, 15)	0
dense_8 (Dense)	(None, 1)	16
=====		
Total params: 496		
Trainable params: 496		
Non-trainable params: 0		

```

1 history = model.fit(x_train, y_train, validation_data = (x_test, y_test), epochs=50, verbose=1)

```

```
1 y_pred = model.predict(x_test)
2 y_pred = (y_pred>0.5)
3 y_pred
```

24/26


```
[False],
[ True],
[False],
[ True],
[False],
[ True],
[ True],
[ True],
[ True],
[ True],
[False],
[False],
[ True],
[False],
[False],
[ True],
[ True],
[False],
[False],
[ True],
[ True],
[False],
[False],
[ True],
[False],
[False],
```

```
1 from sklearn.metrics import confusion_matrix, classification_report, accuracy_score
```

```
1 print(confusion_matrix(y_test, y_pred))
```

```
[[50  0]
 [ 0 50]]
```

```
1 print(classification_report(y_test, y_pred))
```

	precision	recall	f1-score	support
0	1.00	1.00	1.00	50
1	1.00	1.00	1.00	50
accuracy			1.00	100
macro avg	1.00	1.00	1.00	100
weighted avg	1.00	1.00	1.00	100

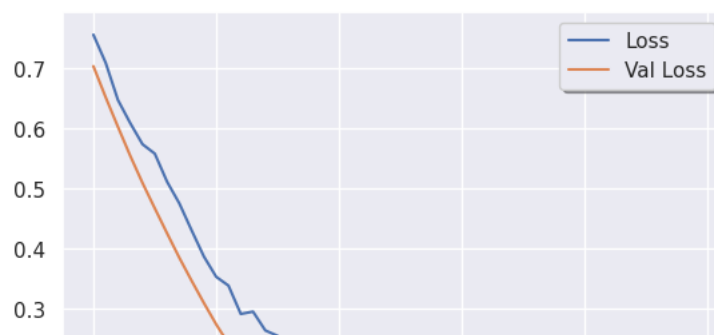
▼ Accuracy is 1 (100%)

```
1 print(accuracy_score(y_test, y_pred))
```

```
1.0
```

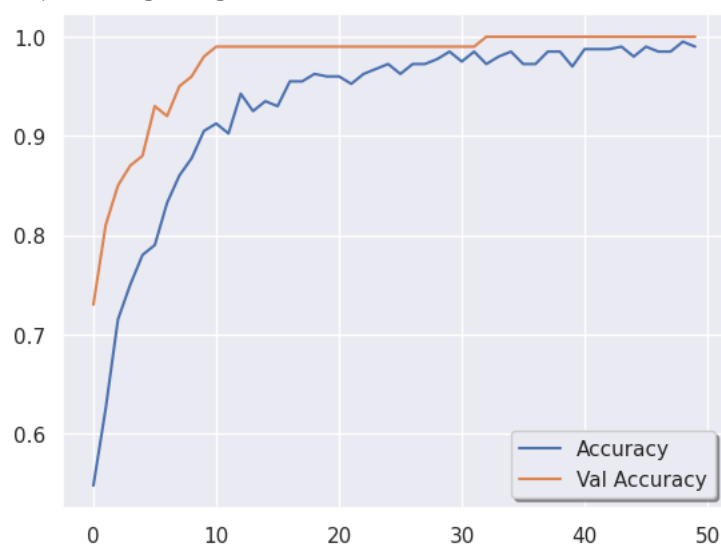
```
1 plt.plot(history.history['loss'])
2 plt.plot(history.history['val_loss'])
3 plt.legend(('Loss', 'Val Loss'),
4           loc='upper right', shadow=True)
```

<matplotlib.legend.Legend at 0x7f7478a99a60>



```
1 plt.plot(history.history['accuracy'], label='test')
2
3 plt.plot(history.history['val_accuracy'])
4 plt.legend(('Accuracy', 'Val Accuracy'),
5           loc='lower right', shadow=True)
```

<matplotlib.legend.Legend at 0x7f7481939340>



Colab paid products - Cancel contracts here

✓ 0s completed at 3:39 PM

